AbstractID:9731Title :Ma chineLearningba sedClinicalRes earch:Theexampleof LungC ancer

## MachineLearnin gbasedC linical Research: Theexampleof  Lung Cancer

P.La mbin[1], C.Dehin g-Oberije[1], L.Persoon [1], S.Yu [2], D.DeRuyssch er[1], R. B.Rao [2], A.Dekker [1].

[1] UniversityM edicalCenter Ma astricht,Depart mentof Ra diationOn cology,MAAS TROclinic,  GROW, TheNetherl ands

[2] SiemensM edicalSolutionsUSA,  Inc., Malvern ,Pe nnsylvania,U SA

**Purpose:** The *hypothesis*o fthislongtermpro   jectist  hata  multicentric based informationsys tembas edon fourmodu les(mul tiparametric interconnectedh ealthcare databases,datamin ing tools,u pdatedmachi ne learning basedpredicti ve algorithmsan duserinter  faces)w illfacilitat ean dacceler ate research inon cology. Weca ll thisap proach"Ma chine LearningBasedClinicalResearch(MLBCR)    ".W e performed apilot projectin  non-small celllungcancer(   NSCLC)p atientsfor w hich clinicalTN M stageis hi ghlyi naccurate fort hepr ediction ofsur vival ofno n-surgicalp atients andalter nativesar e currentlylacking .The  objectives of thisstudy  were todeve lopandvalidateapr   edictionmodelforsur  vivalofNS CLCpatients,  treatedwith (chemo)rad iotherapy, usingclinical  factors.

**Patientsan dMethods:**  Threei nterconnecteddatabasesw eremir roredin  toa  data warehouseus inga disease based,c ohort-specificdatamodel.Thet   hree datasour cesw erea)  electronicmed icalrecords,b ) imaging andDICOM -RTobjects in aRT -PACSandc)  treatmentinf ormationinar  ecordand  verify database. Datafro m403consec  utiveinoperableNSC  LCpati ents,stag eI -IIIB,treated  radically with (chemo)rad iation were selected. In 82 patientsdatafr omb lood samples wereavailable . The2 -norm SupportVe ctor Machines wereusedto  bu ildth epr ognosticmod els. Performanceofthe   modelswas expressedasth eA UC(A reaU ndertheCur  ve)oft heReceiv er OperatingChar acteristic(ROC)  and assessedusi nglea ve-one-out( LOO)cro ss-validation. Thep rognosticmodel,u  sing clinical factors only, wasva lidatedu singtw o external, independentdatasets  with36and65patients,r     espectively. In addition,a

AbstractID:9731 Title :Machine Learning based Clinical Research:The example of Lung Cancer

risk score was calculated and a nomogram,which is in fact a graphical representation of the risk score,was made for practical use.

**Results:** The model , based on 403 patients and using clinical factors,consisted of gender,WHO performance status,forced expiratory volume( $FEV_1$ ),number of positively lymph node stations on PET and gross tumor volume( onPET -CT).The AUC,assessed by LOO cross-validation, was 0. 75(95 %CI 0.70-0.82),while application of the model to the external datasets yielded an AUC of 0.75 and 0.76 respectively. Splitting the MAASTRO cohort into 3 subgroups,based on their risk score,resulted in the identification of a high,medium and low risk group.The 2-year survival was 66%(95%CI 54%-78%) for the lower risk group, 29%( 95%CI 21%-37%) for the medium risk group and 14%(95%CI 5% -23%) for the high risk group. If blood biomarkers were available ,based on the 82 patients the prognostic model consisted of three additional biomarkers factors:OPN, IL8 and CEA.The LOO AUC was 0.83( 95% CI0.76-0.94),which is significantly better than the prognostic model using only clinical factors based on these same 82 patients (AUC0.71,95% CI0.60-0.87). *In conclusion* ,the model,using clinical factors, successfully estimates 2-year survival of NSCLC patients and the performance,assessed internally as well as in two independent datasets,is good. Combining blood biomarkers with clinical factors yielded a significantly better performance than using clinical factors only( AUC:0.83 vs0.71). We concluded that MLBCR is feasible . The bottleneck is the availability of external datasets . Therefore,we need to invest in international standards as well in multicentric approaches allowing to recruit more patients,preferably having had different type of treatments,and to have quick access to external validation datasets.

**Conflict of interest:** This project has been partially funded by Siemens IKM.